



## **Experiment-based detection of service disruption attacks in optical networks using data analytics and unsupervised learning**

Downloaded from: <https://research.chalmers.se>, 2023-05-05 07:14 UTC

Citation for the original published paper (version of record):

Furdek Prekratic, M., Natalino Da Silva, C., Schiano, M. et al (2019). Experiment-based detection of service disruption attacks in optical networks using data analytics and unsupervised learning. Metro and Data Center Optical Networks and Short-Reach Links II; 109460D, 10946. <http://dx.doi.org/10.1117/12.2509613>

N.B. When citing this work, cite the original published paper.

# Experiment-based detection of service disruption attacks in optical networks using data analytics and unsupervised learning

Marija Furdek<sup>a</sup>, Carlos Natalino<sup>a</sup>, Marco Schiano<sup>b</sup>, and Andrea Di Giglio<sup>b</sup>

<sup>a</sup>School of Electrical Engineering and Computer Science, KTH Royal Institute of Technology, Stockholm, Sweden.

<sup>b</sup>Telecom Italia, Turin, Italy.

## ABSTRACT

The paper addresses the detection of malicious attacks targeting service disruption at the optical layer as a key prerequisite for fast and effective attack response and network recovery. We experimentally demonstrate the effects of signal insertion attacks with varying intensity in a real-life scenario. By applying data analytics tools, we analyze the properties of the obtained dataset to determine how the relationships among different optical performance monitoring (OPM) parameters of the signal change in the presence of an attack as opposed to the normal operating conditions. In addition, we evaluate the performance of an unsupervised learning technique, i.e., a clustering algorithm for anomaly detection, which can detect attacks as anomalies without prior knowledge of the attacks. We demonstrate the potential and the challenges of unsupervised learning for attack detection, propose guidelines for attack signature identification needed for the detection of the considered attack methods, and discuss remaining challenges related to optical network security.

**Keywords:** Optical network security, dataset exploration, data analytics, unsupervised learning, anomaly detection.

## 1. INTRODUCTION

As critical communication infrastructure, optical networks have a vital role in safe and dependable transmission of massive amounts of data, supporting essential societal services. However, these networks are inherently vulnerable to a multitude of deliberate attack methods that target service disruption at the physical layer.<sup>1</sup> In spite of their key role, security of optical networks threatened by targeted physical-layer disruptions has largely been overlooked thus far, creating a need for approaches that prepare network operators and enhance network robustness to optical-layer disturbances. This is particularly significant in the current era of ever-increasing scope, intensity and sophistication of attacks. Moreover, in the near future, the existing optical network infrastructure is envisaged to be upgraded with proliferating physical-layer paradigms of quantum key distribution (QKD) and Space Division Multiplexing (SDM), which are highly susceptible to the optical-layer disruptions exacerbated by attacks.

One of the most harmful physical-layer attack methods identified in the literature is the insertion of harmful signals that intensify physical layer impairments and cause degradation of the optical channels. Insertion of harmful signals can take place by obtaining direct access to the patch-panel, or by accessing the optical fiber beyond a secured perimeter, creating a temporary coupler<sup>2</sup> and inserting the signal. Such attacks can strengthen the non-linear effects in fibers and reduce the amplifier gain provided to the useful optical channels. A key prerequisite for fast and effective attack response and network recovery is the timely and accurate attack detection. Identifying the presence of a harmful signal based on such effects would require not only ubiquitous monitoring, but also the capability to recognize the characteristic degradation experienced by the signal and attribute it to a particular attack technique. As OPM equipment<sup>3</sup> is deployed scarcely at strategic network locations due to its restrictive cost, attack diagnostic procedures typically cannot rely on pervasive OPM data available at all

---

M. Furdek: E-mail: marifur@kth.se, Telephone: +46 (0)8 790 4213

More information about the implementation used in this work is available at [GitHub](#).

points in the network, which complicates attack detection. Moreover, the effects caused by harmful signals in the network can differ depending on the signal parameters (e.g., spectrum and power level), as well as the network configuration. Due to a severe lack of experimental data and theoretical models capturing their harmful effects, identifying the signature of attacks remains challenging.

In this paper, we adopt an experimental, data-analytic approach to optical network security diagnostics. First, we experimentally demonstrate the effects of signal insertion attacks with varying intensity in a real-life scenario using a telecom operator testbed. The testbed is equipped with coherent off-the-shelf receivers which provide a rich set of OPM parameters recorded at the destination of the optical signal where it needs to be detected anyway, rather than relying on specialized OPM devices placed at intermediate nodes. By applying data analytic tools, we analyze the properties of the obtained dataset to determine how the relationships among different signal parameters change under an attack scenario as opposed to the normal operating environment. We then apply an unsupervised learning technique, i.e., clustering, to detect attacks as anomalies. We explore how the anomaly detection algorithm can be configured to narrow down what is considered normal operating conditions so as to increase the accuracy of identifying the attacks without prior knowledge of their effects. The proposed guidelines on attack signature identification contribute to the detection of the considered attack methods.

This paper is organized as follows. Sec. 2 presents an overview of related work from the literature. Sec. 3 presents the experimental setup used to acquire the data, whose properties are then investigated in Sec. 4. Sec. 5 provides concluding remarks on the presented approaches and on remaining challenges in optical network security management.

## 2. LITERATURE REVIEW

A general overview of security vulnerabilities associated to the physical layer of optical networks, primarily in the context of attacks aimed at service disruption, can be found in Ref. 1. The authors in Ref. 4 conducted an evaluation of the damage from out-of-band jamming on co-propagating optical channels in simple Wavelength Division Multiplexing (WDM) networks via simulations. Sensitivity of novel optical transmission paradigms, such as QKD and SDM, to physical-layer impairments and disruptions was studied in Refs. 5–7. In particular, the work in Ref. 5 evaluated the impact of the co-propagating classical channel launch power on the QKD system’s secret key rate, indicating their strong inverse proportionality. The works in Refs. 6, 7 studied the impact of excessive channel power on the SDM system performance, indicating that jamming attacks can reduce transmission reach and strongly impact channels that use higher-order modulations formats in adjacent cores of a multi-core fiber (MCF).

Detection of attacks requires monitoring of optical channels’ performance and the analysis of the collected data. Monitoring of optical parameters from the network and combining it with monitoring data from other layers is enabled by Software-Defined Optical Networking (SDON).<sup>8</sup> The advances in coherent receivers with digital signal processing (DSP) capabilities provide an extensive OPM dataset for each received channel at their destination, and enable high-performance monitoring with short monitoring cycles.<sup>9,10</sup> However, this generates a large amount of data which needs to be processed by the SDON controller. For this reason, several works in the literature rely on the application of machine learning (ML) for monitoring and managing the performance of optical networks in a predictive and a reactive manner.<sup>11–13</sup>

The application of ML to assist in complex optical network management decisions has been shown to bring very promising advantages.<sup>14,15</sup> In dynamic network operation, ML algorithms have been used to model complex optical network operations such as service provisioning,<sup>16–18</sup> admission control,<sup>19</sup> and traffic prediction.<sup>18,20</sup> These ML algorithms have shown potential benefits towards supporting autonomous monitoring and operation of optical networks. For instance, ML-assisted failure prediction on optical networks can substantially increase the availability while reducing the need to post-failure reconfigurations.<sup>21</sup>

In Ref. 22, we have applied ML techniques to the field of optical network security for the first time. The results therein indicated a strong potential of Artificial Neural Networks (ANNs) to detect out-of-band jamming signals of different intensities (i.e., power levels) with an average accuracy of 93%. In this paper, we perform

a deeper analysis of the experimental dataset obtained by performing out-of-band jamming attacks in a field-deployed operator testbed, with the goal to identify the main trends in OPM parameter correlations that could serve as an indicator of attacks in wider contexts. Moreover, we apply unsupervised learning to understand how, unlike in Ref. 22, attacks can be detected without prior knowledge of their characteristics.

### 3. EXPERIMENTAL SETUP FOR DATASET ACQUISITION

Field-deployed experimental setup depicted in Fig. 1 was used to perform the experiments. The setup comprises a Coriant Groove G30 coherent transponder, two Flexgrid ROADMs and a field-deployed optical line system with 4 amplification sections and 280 km total length. The optical channel under test carries 200 Gbps with nominal center frequency at 193.1 THz. 6 CW channels are used in order to simulate realistic loading conditions. All channels have the same launch power.

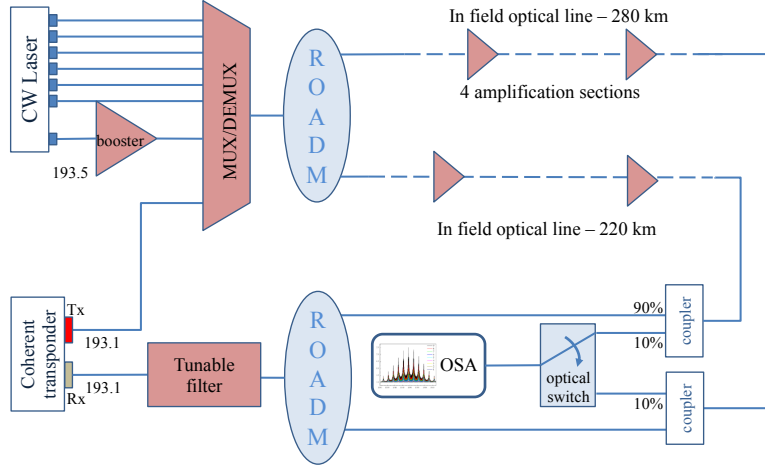


Figure 1. Setup used in the experiments.

In the implemented attack scenario, a jamming signal was inserted in the optical multiplexer at 193.5 THz, and its power levels were varied. At the beginning of the experiment, the jamming signal was switched off, simulating normal operating conditions. The power of the jamming signal was then increased to 0 dB, 3 dB and 6 dB (with respect to the other channels), corresponding to the light, moderate and strong attack condition, respectively. Table 1 shows the OPM parameters collected during our experiments and their respective units. We gather a balanced dataset comprising the same number of samples for each attack condition, where the OPM parameters are collected every minute by an automated procedure.

### 4. EXPLORATORY DATASET ANALYSIS

In this section, we perform a detailed analysis of the obtained dataset properties, and investigate the applicability of several data analytic approaches to support autonomous identification of attacks. We used Python 3, pandas<sup>23</sup> and Scikit-learn<sup>24</sup> for the implementation of all the procedures described in this work.\*

#### 4.1 Dataset characteristics

The monitored values of the parameters are shown in Fig. 2a, while Fig. 2b depicts their standardized values. The standardization technique scales the values by computing their distance from the average in terms of standard deviation, in such a way that the correlation between the OPM parameters is maintained. As can be seen in Fig. 2a, the collected OPM parameters can assume values of very different orders of magnitude, which highlights the importance of standardizing them before further analysis, as ML algorithms usually perform better with standardized data.

\*More information about the implementation used in this work is available at [GitHub](#).

Table 1. Parameters collected from the testbed using optical performance monitoring (OPM).<sup>22</sup>

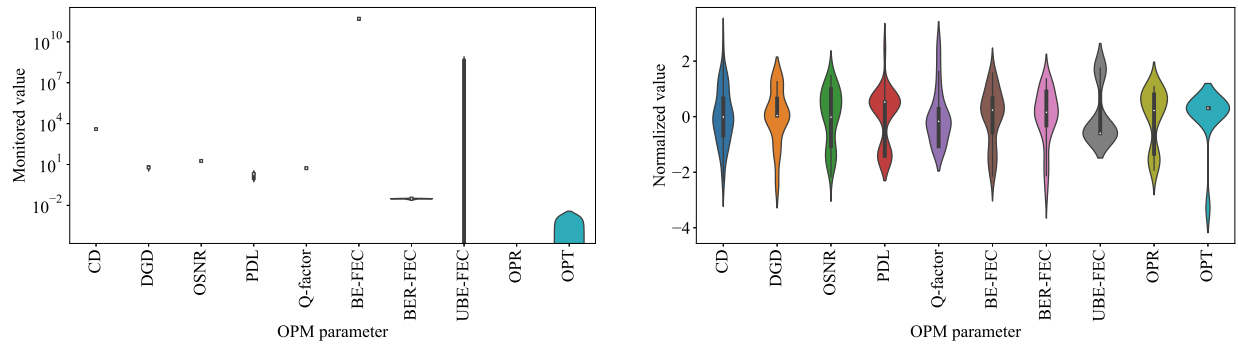
Acronym	Unit	Description
CD	ps/nm	Chromatic Dispersion
DGD	ps	Differential Group Delay
OSNR	dB	Optical Signal to Noise Ratio
PDL	dB	Polarization Dependent Loss
Q-factor	dB	Q factor
BE-FEC	Bits	Block Errors before FEC
BER-FEC	Bps	Bit Error Rate before FEC
UBE-FEC	Blocks	Uncorrected Block
BER-POST-FEC	Bps	Bit Error Rate after FEC
OPR	dBm	Optical Power Received
OPT	dBm	Optical Power Transmitted
OFT	MHz	Optical Frequency Transmitted
OFR	MHz	Optical Frequency Received

For all these parameters except BE-FEC and UBE-FEC, the system provides the maximum, minimum and average value in the observation interval.

A statistical procedure which can be used to analyze the experimentally gathered data and identify the attack scenario is Principal Component Analysis (PCA). PCA transforms the data to map it into distinct, orthogonal principal components, each of which has the largest possible variance. Fig. 3 shows the results of PCA applied to the standardized dataset. Applying PCA indicates the level of difficulty in distinguishing between the normal operating conditions and the different attack scenarios. Fig. 3 shows that achieving clear separation of light and moderate attacks from the normal operating condition can be challenging. Therefore, a deeper analysis is needed in order to obtain clearer understanding of the attack properties and effects.

## 4.2 Correlation between OPM parameters

We begin the dataset exploration by investigating the trends in correlation between different OPM parameters. Correlation refers to the statistical relation between pairs of parameters, which can range from -1 (the strongest



(a) Monitored values (with y axis in log scale)

(b) Standardized values

Figure 2. Statistical representation of the collected dataset for the monitored and standardized values. The outer shape represents the distribution of values across y. Inner boxes show 50% and 95% probability, respectively. White dots represent the median average. Note that the standardization makes the features have an average near to zero.

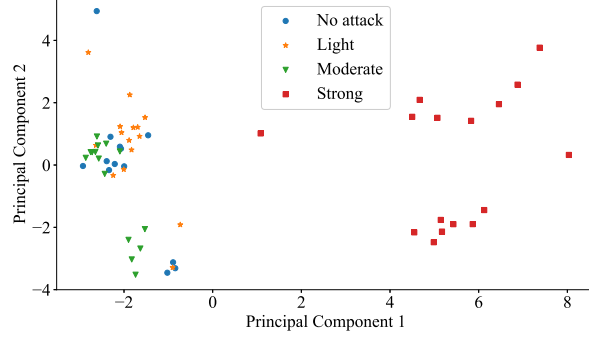


Figure 3. PCA for the standardized dataset.

possible disagreement) to 1 (the strongest possible agreement).

Fig. 4 shows the trends in correlation between OPM parameters under different operating conditions. To ease tracking, the values are color-coded according to the heat map shown on the right hand side. Fields with no shown numerical value are *NaN*. Fig. 4a shows the OPM parameter correlation under normal, attack-free operating condition. Observations that can be made from the attack-free samples include a strong negative correlation between optical power received (OPR) and pre-FEC bit error rate (BER-FEC), and a strong positive correlation between OPR and Q-factor, which is in agreement with the theoretical relationships of these parameters (i.e., higher received signal power under normal conditions results in higher Optical Signal-to-Noise Ratio (OSNR) and decreases the error rates).

Figs. 4b–4d depict the differences in the correlation of OPM parameters under light, moderate, and strong attack, respectively, compared to the no-attack scenario as baseline. The difference  $\Delta_{corr}$  between parameter correlation in a particular attack scenario, denoted as  $corr_{attack}$ , and the normal condition, denoted as  $corr_{noattack}$ , is computed as:

$$\Delta_{corr} = corr_{attack} - corr_{noattack}. \quad (1)$$

As can be seen in Figs. 4b–4d, the presence of a harmful signal reverses the correlation between some OPM parameters compared to the no-attack scenario. A prominent example of this is the aforementioned correlation between OPR and BER-FEC, which becomes positive in the presence of an attack. This can be explained by the fact that the noise exacerbated by the harmful signal raises the received power of the optical channel under test. Another example is the correlation between OPR and Q-factor, whose value takes a significant turn towards negativity in the light and strong attack scenario.

### 4.3 Correlation-based attack detection

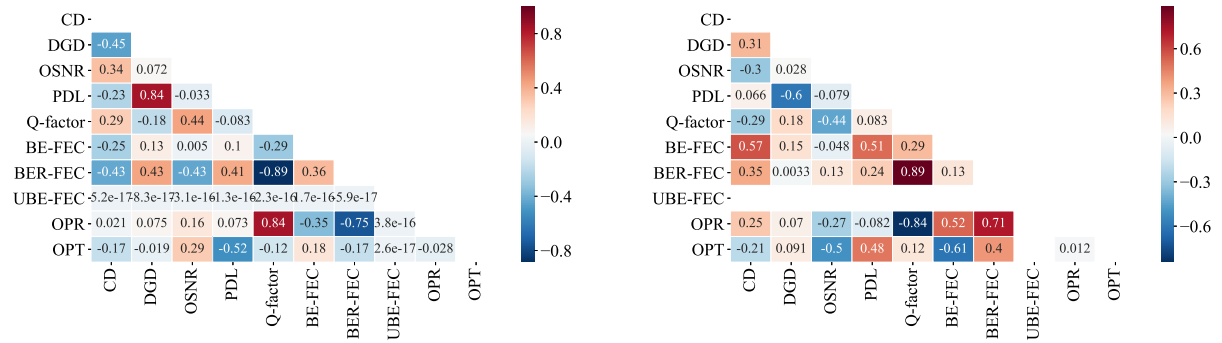
Based on the knowledge about the correlation between OPM parameters obtained in the previous section, we develop a correlation-based attack detection approach. This approach mimics the reasoning of a specialist developing an attack detection procedure based only on prior knowledge and correlation. Correlation can only be used if we consider a specific time window associated to the number of most recent monitoring samples used to compute correlation (denoted by  $\tau$ ). For instance, if we consider  $\tau = 20$ , the last 20 monitoring samples are used to compute correlation between the OPM parameters. The higher is the number of samples  $\tau$ , the more stable is the analysis, and it becomes less likely that the usual variations in the OPM parameters will cause large variations in the correlation. However, higher values of  $\tau$  incur a greater delay in detecting the correlation changes observed during attacks because more monitoring samples need to be collected during the attack to cause a change in the correlation values. In addition, this method requires the setting of a threshold for the correlation values which trigger attack detection alarms.

Fig. 5 illustrates correlation-based attack detection for our experimental dataset considering different values of  $\tau$ . For the sake of clarity, in this scenario, only the normal operation (i.e., no attack) condition and the

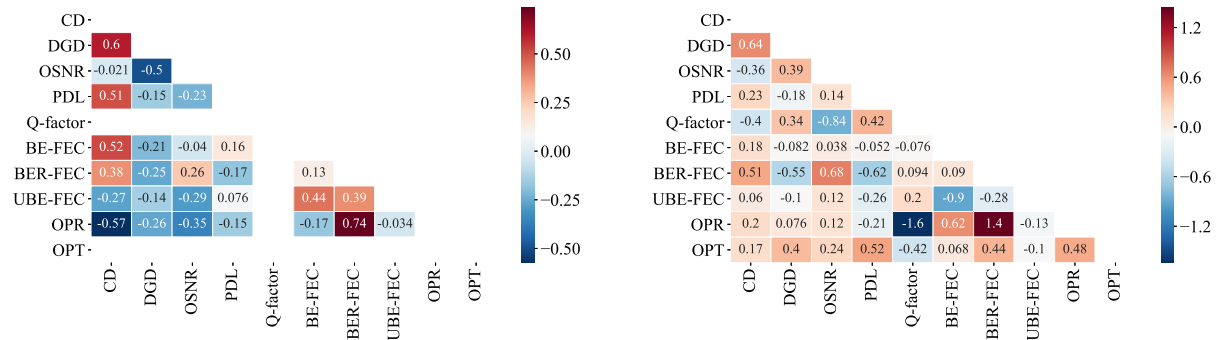
strong attack condition are considered. The red-dotted vertical line represents the transition from the no-attack to the strong attack condition. As we consider 1-minute resolution of OPM parameter collection, the  $x$ -axis values correspond to the elapsed time and the number of collected samples. When  $\tau = 20$ , (Fig. 5a), around the 5<sup>th</sup> sample, the two correlation values sharply change towards zero. This is caused by small variations in the OPM parameters rather than an attack, and illustrates how a too low number of considered samples can disrupt correlation and cause a false alarm. Considering the same data with  $\tau = 30$  (Fig. 5b), this sudden change in the correlation trends under no attack condition is eliminated.

The benefit of considering a lower value of  $\tau$  is observed in the time it takes to detect the attack. Setting  $\tau = 20$  introduces a change in the correlation trends at the second monitoring sample upon the attack. For the 1-minute monitoring window considered, the attack is detected in 2 minutes. With  $\tau$  set to 30, the slope of the considered correlation curves is more gradual, leading to longer detection time. Here, a clear correlation change is observed after four monitoring samples, i.e., in 4 minutes. The time to detect an attack could be shortened by setting a tight correlation threshold to trigger an alarm, but this can in turn produce a higher number of false positives. On the other hand, supervised learning approaches which require prior knowledge of the attacks, such as the ones presented in Ref. 22, can detect attacks within a single monitoring cycle. This indicates a clear trade-off between prior knowledge and the time needed to detect an attack.

Although correlation-based attack detection presents promising results, it requires a deep understanding of the physical consequences of attacks. Moreover, it is challenging to define a set of parameters that can help obtain low false positive and false negative rates. In the following section, we explore other algorithms that can potentially detect attacks without prior knowledge of their signatures.



(a) Parameter correlation in normal, no-attack condition. (b) Correlation difference between no-attack and light attack condition.



(c) Correlation difference between no-attack and moderate attack condition. (d) Correlation difference between no-attack to strong attack condition.

Figure 4. Correlation between the OPM parameters in different attack scenarios.



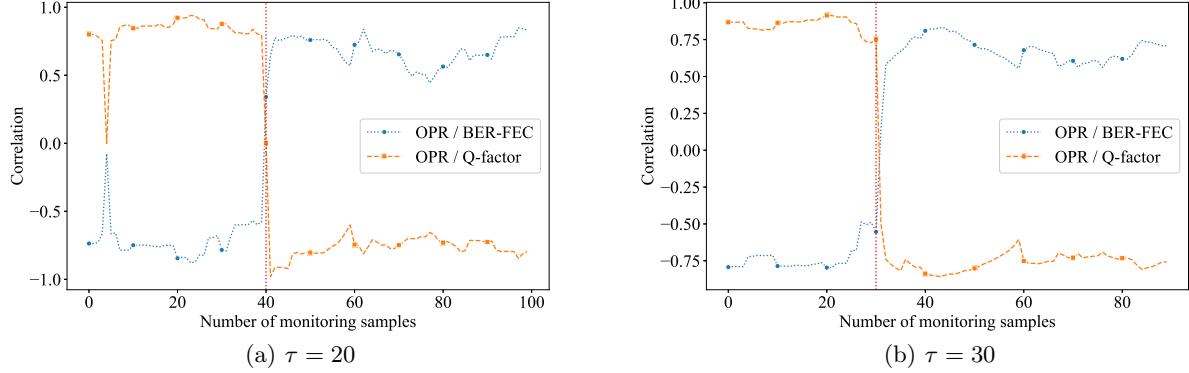


Figure 5. Variation of correlation over time considering the last 20 or 30 monitoring samples. The vertical red dotted line represents the start of the attack.

#### 4.4 Clustering-based attack detection

The results presented in the previous section rely on a deep knowledge of the physical effects of jamming signal insertion attacks. However, with the increasing complexity of optical networks and the expanding threat landscape, it is highly desirable for a security assessment system to detect anomalies without prior knowledge. To this end, we employ the density-based spatial clustering of applications with noise (DBSCAN) algorithm,<sup>25</sup> which has been successfully applied to anomaly detection in optical networks.<sup>13</sup> The DBSCAN has two parameters, denoted as  $\epsilon$  and  $MinPts$  (see Ref. 25 for more details), that must be set in order to properly detect outliers, which are the attacks to be identified. The  $\epsilon$  parameter sets the maximum distance between samples to be considered neighbors, while  $MinPts$  defines the minimum number of neighbors for a sample to be considered a core node. In our work, we consider the Euclidean distance to compute the distance between samples.

Many works in the literature, such as Ref. 13, assume the existence of some anomalies in the dataset already at the beginning of the anomaly detection setup. We take a different approach by starting with a system where only the normal operation conditions are known, i.e., there are no known occurrences of attacks so far. This is realistic for real-world networks experiencing previously unseen techniques of attacks. The main objective here is to set the parameters in such a way that it becomes very hard to identify an anomaly as no attack condition. When only normal operating conditions are known, DBSCAN should be configured in such a way that no false positives are observed. In this case, false negatives are not possible since the dataset only has samples from normal, i.e., no-attack condition.

Fig. 6 shows the percentage of false positives obtained by DBSCAN for different configurations of  $\epsilon$  and

1	0.00	0.00	0.00	0.00	0.00	0.00
3	3.33	3.33	3.33	3.33	0.00	0.00
5	21.67	21.67	10.00	10.00	3.33	0.00
8	100.00	68.33	46.67	46.67	21.67	0.00
10	100.00	68.33	46.67	46.67	21.67	13.33
12	100.00	68.33	46.67	46.67	40.00	13.33
15	100.00	68.33	68.33	68.33	40.00	13.33
20	100.00	100.00	100.00	100.00	40.00	40.00
	0.1	0.5	1.0	1.0	2.0	3.0
	$\epsilon$					

Figure 6. Percentage of no attack samples classified as attacks (i.e., false positives) obtained by DBSCAN with different parameter settings for the dataset with only attack-free samples.



*MinPts*. In this case, low *MinPts* and high  $\epsilon$  values yield the lowest false positive rates. Considering a real-world scenario, one of the settings with zero false positives would be chosen for the attack detection algorithm. It is important to mention, however, that these parameters might need further tuning if the dataset properties, e.g., the number of samples, change.

Once the performance of DBSCAN in the presence of only attack-free samples has been analyzed, we move our attention to the performance of the DBSCAN in the presence of malicious attacks. We consider that the attack samples are present in a 1:12 proportion with respect to the no-attack condition, assuming that in real-world conditions attacks would occur rarely. We assess the performance in terms of false positives and false negatives for each of the three attack scenarios. This allows us to individually scrutinize the performance for more challenging attacks such as the light one.

Table 2 presents the percentage of false positives ( $f_p$ ) and false negatives ( $f_n$ ), considering a dataset composed of the no-attack samples and the ones from the specific attack scenario identified in the top row. As shown in Fig. 6, for *MinPts*=1 and for *epsilon*=4.0, there are no false positives. For light and moderate attacks, it is very

Table 2. False positive and false negative rates ( $f_p$ ,  $f_n$ ) obtained by DBSCAN in different attack scenarios and for different parameter settings.

	$\epsilon$						
<i>MinPts</i>	0.1	0.5	1.0	1.0	2.0	3.0	4.0
	<b>Light attack</b>						
1	0.00, 7.69	0.00, 7.69	0.00, 7.69	0.00, 7.69	0.00, 7.69	0.00, 7.69	0.00, 7.69
3	3.08, 0.00	3.08, 0.00	3.08, 3.08	3.08, 3.08	3.08, 7.69	0.00, 7.69	0.00, 7.69
5	21.54, 0.00	15.38, 0.00	15.38, 3.08	15.38, 3.08	3.08, 7.69	0.00, 7.69	0.00, 7.69
8	78.46, 0.00	64.62, 0.00	38.46, 3.08	38.46, 3.08	10.77, 6.15	0.00, 7.69	0.00, 7.69
10	92.31, 0.00	64.62, 0.00	38.46, 3.08	38.46, 3.08	23.08, 6.15	23.08, 6.15	0.00, 7.69
12	92.31, 0.00	64.62, 0.00	38.46, 3.08	38.46, 3.08	23.08, 6.15	23.08, 6.15	0.00, 7.69
15	92.31, 0.00	64.62, 0.00	38.46, 3.08	38.46, 3.08	23.08, 6.15	23.08, 6.15	0.00, 7.69
	<b>Moderate attack</b>						
1	0.00, 7.69	0.00, 7.69	0.00, 7.69	0.00, 7.69	0.00, 7.69	0.00, 7.69	0.00, 7.69
3	3.08, 0.00	3.08, 0.00	3.08, 0.00	3.08, 0.00	3.08, 7.69	0.00, 7.69	0.00, 7.69
5	21.54, 0.00	15.38, 0.00	15.38, 0.00	15.38, 0.00	3.08, 7.69	0.00, 7.69	0.00, 7.69
8	78.46, 0.00	64.62, 0.00	38.46, 0.00	38.46, 0.00	10.77, 7.69	10.77, 7.69	0.00, 7.69
10	92.31, 0.00	64.62, 0.00	38.46, 0.00	38.46, 0.00	23.08, 7.69	23.08, 7.69	0.00, 7.69
12	92.31, 0.00	64.62, 0.00	38.46, 0.00	38.46, 0.00	23.08, 7.69	23.08, 7.69	0.00, 7.69
15	92.31, 0.00	64.62, 0.00	38.46, 0.00	38.46, 0.00	23.08, 4.62	23.08, 7.69	0.00, 7.69
	<b>Strong attack</b>						
1	0.00, 7.69	0.00, 7.69	0.00, 7.69	0.00, 7.69	0.00, 7.69	0.00, 7.69	0.00, 7.69
3	3.08, 0.00	3.08, 0.00	3.08, 0.00	3.08, 0.00	3.08, 0.00	0.00, 0.00	0.00, 4.62
5	21.54, 0.00	15.38, 0.00	15.38, 0.00	15.38, 0.00	3.08, 0.00	0.00, 0.00	0.00, 4.62
8	78.46, 0.00	64.62, 0.00	38.46, 0.00	38.46, 0.00	10.77, 0.00	10.77, 0.00	0.00, 4.62
10	92.31, 0.00	64.62, 0.00	38.46, 0.00	38.46, 0.00	23.08, 0.00	23.08, 0.00	0.00, 4.62
12	92.31, 0.00	64.62, 0.00	38.46, 0.00	38.46, 0.00	23.08, 0.00	23.08, 0.00	0.00, 4.62
15	92.31, 0.00	64.62, 0.00	38.46, 0.00	38.46, 0.00	23.08, 0.00	23.08, 0.00	0.00, 1.54

challenging to achieve a low false negative rate without degrading the false positive rate. For instance, in order to achieve zero false positive rate, the lowest false negative for light and moderate attacks is 7.69% (denoted with red in the table). However, in attack detection, false negatives are more harmful than false positives, as they lead to attacks being undetected. To address this, it is possible to configure the algorithm in such a way that no false negatives are observed, at the expense of a 3.08% false positive rate (denoted with blue in the table). For the strong attack case, selecting  $MinPts=3$  or  $MinPts=5$  and  $\epsilon=3.0$  leads to zero false positives and zero false negatives (denoted with green in the table).

By analyzing Fig. 6 and Table 2, it is possible to extract some interesting insights. For instance, configuring the algorithm with the best results from Fig. 6, i.e.,  $\epsilon = 4.0$  already ensures no false positives and a maximum of 7.69% of false negatives. This means that such an algorithm can be helpful even if possible threats are completely unknown. Moreover, increasing  $MinPts$  makes it more difficult to have attack samples classified as no attack, and balancing it with a good  $\epsilon$  configuration can further enhance the performance of the algorithm. Finally, the configuration can be fine-tuned periodically to cope with new threats that may appear.

## 5. CONCLUSION

In this paper, we performed an analysis of a dataset collected from real-world testbed in the presence of jamming attacks exploring different ways of detecting attacks. First, we analyzed which OPM parameters could be used for a correlation-based attack detection. Then, we used unsupervised learning to detect attacks without prior knowledge.

We showed that it is possible to detect attacks by analyzing the correlation between key OPM parameters, but this approach may lead to attacks being undetected for a relatively long time before triggering the attack alarm. We also showed that it is possible to configure an anomaly detection approach in such a way that unknown attacks can be detected in 92% of the cases (i.e., at a false negative rate of 7.69%). Having prior knowledge of attacks can improve the performance and decrease the undetected attack occurrences to 0% for moderate and strong attacks, by sacrificing the false positive rate of 3%. This shows that supervised learning approaches, which can detect attacks with a single OPM reading of the attack, are still the most accurate and quickest attack detection method, with the drawback of requiring prior knowledge of the attacks. Further developing unsupervised learning methods is necessary to achieve a truly autonomous attack detection.

There are still many open issues related to the security of optical networks. Specifically for the case of analytics-based attack detection, exploring mathematical properties other than the correlation explored in this work might present interesting possibilities. Regarding the application of ML to the optical network security, feature engineering could be performed to understand whether the unsupervised learning algorithm can be disturbed/improved or not, upon presenting only a part of the collected OPM parameters.

## ACKNOWLEDGMENTS

We gratefully acknowledge Coriant for providing the Groove G30 transponder. This article is based upon work from COST Action 15127 RECODIS and Celtic-Plus project SENDATE-EXTEND.

## REFERENCES

- [1] Skorin-Kapov, N. et al., “Physical-layer security in evolving optical networks,” *IEEE Commun. Mag.* **54**(8), 110–117 (2016).
- [2] Uematsu, T. et al., “Design of a temporary optical coupler using fiber bending for traffic monitoring,” *IEEE Photonics J.* **9**, 1–13 (Dec 2017).
- [3] Dong, Z. et al., “Optical performance monitoring: A review of current and future technologies [invited],” *IEEE/OSA J. Lightwave Techn.* **34**, 252–543 (Jan 2016).
- [4] Peng, T. et al., “Propagation of all-optical crosstalk attack in transparent optical networks,” *Opt. Eng.* **50**, 085002.1–3 (Aug 2011).
- [5] Ou, Y. et al., “Field-trial of machine learning-assisted quantum key distribution (QKD) networking with SDN,” in [*Proc. of ECOC*], Mo3D.3 (Sept 2018).

- [6] Perelló, J. et al., “Flex-grid/SDM backbone network design with inter-core XT-limited transmission reach,” *IEEE/OSA J. Opt. Commun. and Netw.* **8**, 540–552 (Aug 2016).
- [7] Goścień, R. et al., “Impact of high-power jamming attacks on SDM networks,” in [*Proc. of ONDM*], 77–81 (May 2018).
- [8] Thyagaturu, A. S. et al., “Software defined optical networks (SDONs): A comprehensive survey,” *IEEE Commun. Surveys Tuts.* **18**, 2738–2786 (Fourthquarter 2016).
- [9] Vela, A. P., Ruiz, M., and Velasco, L., “Distributing data analytics for efficient multiple traffic anomalies detection,” *Computer Communications* **107**, 1 – 12 (2017).
- [10] Paolucci, F. et al., “Network telemetry streaming services in sdn-based disaggregated optical networks,” *IEEE/OSA J. Lightw. Technol.* **36**, 3142–3149 (Aug 2018).
- [11] Shahkarami, S. et al., “Machine-learning-based soft-failure detection and identification in optical networks,” in [*Proc. of OFC*], M3A.5 (2018).
- [12] Vela, A. P. et al., “Soft failure localization during commissioning testing and lightpath operation,” *IEEE/OSA J. Opt. Commun. Netw.* **10**, A27–A36 (Jan 2018).
- [13] Chen, X. et al., “On real-time and self-taught anomaly detection in optical networks using hybrid unsupervised/supervised learning,” in [*Proc. of ECOC*], We1D.4 (Sept 2018).
- [14] Mata, J. et al., “Artificial intelligence (AI) methods in optical networks: A comprehensive survey,” *Optical Switching and Networking* **28**, 43 – 57 (2018).
- [15] Musumeci, F. et al., “An overview on application of machine learning techniques in optical networks,” *IEEE Commun. Surveys Tuts.* **99**, 1–1 (2018).
- [16] Chen, X. et al., “Deep-RMSA: A deep-reinforcement-learning routing, modulation and spectrum assignment agent for elastic optical networks,” in [*Proc. of OFC*], W4F.2 (March 2018).
- [17] Natalino, C. et al., “Machine-learning-based routing of qos-constrained connectivity services in optical transport networks,” in [*Proc. of OSA Advanced Photonics*], NeW3F.5 (2018).
- [18] Natalino, C. et al., “Machine learning aided orchestration in multi-tenant networks,” in [*Proc. of IEEE Photonics Society Summer Topical Meeting Series (SUM)*], 125–126 (July 2018).
- [19] Raza, M. R. et al., “A slice admission policy based on reinforcement learning for a 5g flexible ran,” in [*Proc. of ECOC*], Mo3D.5 (Sept 2018).
- [20] Guo, J. et al., “When deep learning meets inter-datacenter optical network management: Advantages and vulnerabilities,” *Journal of Lightwave Technology* **36**, 4761–4773 (Oct 2018).
- [21] Natalino, C. et al., “A proactive restoration strategy for optical cloud networks based on failure predictions,” in [*Proc. of ICTON*], 1–5 (July 2018).
- [22] Natalino, C. et al., “Field demonstration of machine-learning-aided detection and identification of jamming attacks in optical networks,” in [*Proc. of ECOC*], We2.58 (Sept 2018).
- [23] McKinney, W., “Data structures for statistical computing in python,” in [*Proceedings of the 9th Python in Science Conference*], van der Walt, S. and Millman, J., eds., 51 – 56 (2010).
- [24] Pedregosa, F. et al., “Scikit-learn: Machine learning in Python,” *Journal of Machine Learning Research* **12**, 2825–2830 (2011).
- [25] Ester, M. et al., “A density-based algorithm for discovering clusters in large spatial databases with noise,” in [*Kdd*], **96**(34), 226–231 (1996).